

Course Content

Two-sample analysis

Approximations of expectation and variance

Transformations

Regression Analysis

Two-level experiments

Analysis of variance

Statistical Process Control

Chi Square tests for distributions, independence and homogeneity

Non Parametric Statistics

Two sample tests

Random sample based methods (Z-test and t-test) should not be used on correlated data.

Historical data may help but need to be representative.

Approximation of expectation and variance and transformations

Approximation of mean and variance. $Y = f(X_1, \dots, X_n)$. For independent variables and small variances:

$$E[Y] \approx f(\mu_1, \dots, \mu_n)$$

$$Var[Y] \approx \sum_{i=1}^n \left(\frac{\partial f(\mu_1, \dots, \mu_n)}{\partial x_i} \right)^2 Var[X_i]$$

Transformation of variables to obtain constant variance.

$$Y = g(X) \approx g(\mu) + g'(\mu)(X - \mu)$$

$$Var[Y] \approx (g'(\mu))^2 Var(X) = k$$

$$\sigma_X \propto \mu \Rightarrow g(X) = \ln(X)$$

$$\sigma_X \propto \mu^2 \Rightarrow g(X) = \frac{1}{X}$$

$$\sigma_X \propto \sqrt{\mu} \Rightarrow g(X) = \sqrt{X}$$

Simple linear regression

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \begin{cases} E[\varepsilon_i] = 0, \text{Var}[\varepsilon_i] = \sigma^2 \\ \text{independent} \end{cases}$$

Method of least squares: Minimixe $Q(b_0, b_1) = \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2$ with respect to b_0 and b_1 . This gives:

$$\begin{aligned} b_1 &= \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}, & \hat{\beta}_1 &= \frac{\sum_{i=1}^n (x_i - \bar{x}) Y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ b_0 &= \bar{y} - b_1 \bar{x}, & \hat{\beta}_0 &= \bar{Y} - \hat{\beta}_1 \bar{x} \end{aligned}$$

Regression line: $\hat{y} = b_0 + b_1 x$

Residuals: $y_i - \hat{y}_i = y_i - b_0 - b_1 x_i$

Inference in linear regression

$$T_{\beta_1} = \frac{\hat{\beta}_1 - \beta_1}{S} \sim t_{n-2}$$

$$\sqrt{\frac{S_{xx}}{\sum_{i=1}^n x_{1i}^2}}$$

$$T_{\beta_0} = \frac{\hat{\beta}_0 - \beta_0}{S \sqrt{\frac{\sum_{i=1}^n x_{1i}^2}{n S_{xx}}}} \sim t_{n-2}$$

$$\text{where } S_{xx} = \sum_{i=1}^n (x_{1i} - \bar{x}_1)^2$$

Confidence interval for expected value at x_0 :

$$\hat{y}_0 \pm t_{\frac{\alpha}{2}, n-2} s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x}_1)^2}{S_{xx}}}$$

Prediction interval at x_0

$$\hat{y}_0 \pm t_{\frac{\alpha}{2}, n-2} s \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x}_1)^2}{S_{xx}}}$$

Multiple Linear regression

$$Y_i = \beta_0 + \beta_1 x_{1i} + \cdots + \beta_k x_{ki} + \varepsilon_i$$

$\varepsilon_i \sim N(0, \sigma^2)$ and independent, $i = 1, 2, \dots, n$

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y}$$

$$Cov(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}' \mathbf{X})^{-1}$$

$$\hat{\mathbf{Y}} = \mathbf{X} \hat{\boldsymbol{\beta}} \quad \text{and fitted model: } \hat{\mathbf{y}} = \mathbf{X} \mathbf{b}$$

$$\hat{\sigma}^2 = S^2 = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n - k - 1}$$

Partitioning of variation

$$\sum_{i=1}^n (y_i - \bar{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y}_i)^2$$

$$\text{or } SS_T = SS_E + SS_R$$

Inference Multiple regression

Confidence interval for the mean in x_0 , Multiple regression

$$\hat{y}_0 \pm t_{\frac{\alpha}{2}, n-k-1} s \sqrt{\mathbf{x}_0' (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x}_0}$$

Prediction interval for a new observation in x_0 , Multiple regression

$$\hat{y}_0 \pm t_{\frac{\alpha}{2}, n-k-1} s \sqrt{1 + \mathbf{x}_0' (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x}_0}$$

Partial F-tests

$H_0 : \beta_{i+1} = \beta_{i+2} = \dots = \beta_k = 0$ H_1 : at least one is different from zero.

$$F = \frac{\frac{SS_R(\beta_1, \dots, \beta_k) - SS_R(\beta_1, \dots, \beta_i)}{k-i}}{\frac{SS_E(\beta_1, \dots, \beta_k)}{n-k-1}} \sim F_{k-i, n-k-1}$$

Forward selection

$$\max F_{in} = \frac{\frac{SS_R(\beta_1, \dots, \beta_{i+1}) - SS_R(\beta_1, \dots, \beta_i)}{1}}{\frac{SS_E(\beta_1, \dots, \beta_{i+1})}{n-(i+2)}} \sim F_{1, n-(i+2)}$$

Backward elimination

$$\min F_{in} = \frac{\frac{SS_R(\beta_1, \dots, \beta_{i+1}) - SS_R(\beta_1, \dots, \beta_i)}{1}}{\frac{SS_E(\beta_1, \dots, \beta_{i+1})}{n-(i+2)}} \sim F_{1, n-(i+2)}$$

Residuals in multiple linear regression

$$\hat{\boldsymbol{\varepsilon}} = (\mathbf{I} - \mathbf{H})\boldsymbol{\varepsilon} \text{ where } \mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$$

Studentized residuals:

$$r_i = \frac{e_i}{s\sqrt{1-h_{ii}}}$$

Check of

1. Outliers
2. Heterogeneity in variance
3. Misspecified model
4. Normal distribution

Quality of fit

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$$R_{adj}^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \cdot \frac{n-1}{n-k-1}$$

Quality of prediction

$$R_{pred}^2 = 1 - \frac{\sum_{i=1}^n \frac{(y_i - \hat{y}_i)^2}{(1-h_{ii})^2}}{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

Orthogonal design matrix

X orthogonal

$$\hat{\beta}_i = (\mathbf{x}_i' \mathbf{x}_i)^{-1} \mathbf{x}_i' \mathbf{Y}$$

$$SS_R(\beta_0, \dots, \beta_k) = \sum_{i=1}^k R(\beta_i)$$

Two level designs

Forsøksnr.	Temperatur	Mengde karbon i %	Prosentvis fjærer utan sprekkdanning
1	1450 F	0.5	67
2	1600 F	0.5	79
3	1450 F	0.7	61

4	1600F	0.7	75
---	-------	-----	----

$$Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_{12} x_{1i} x_{2i}$$

$$X = \begin{bmatrix} 1 & 1450 & 0.5 & 725 \\ 1 & 1600 & 0.5 & 1015 \\ 1 & 1450 & 0.7 & 1120 \\ 1 & 1600 & 0.7 & 7200 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 67 \\ 79 \\ 61 \\ 75 \end{bmatrix}, \quad \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ b_{12} \end{bmatrix} = \begin{bmatrix} 14.3 \\ 0.047 \\ -126.7 \\ 0.067 \end{bmatrix}$$

$$x_A = \frac{A - 1525}{75}, \quad x_B = \frac{B - 0.6}{0.1}$$

$$X = \begin{bmatrix} 1 & -1 & -1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & 1 & 1 \end{bmatrix}, \quad \begin{bmatrix} b_0^c \\ b_1^c \\ b_2^c \\ b_{12}^c \end{bmatrix} = \begin{bmatrix} 78.5 \\ 6.5 \\ -2.5 \\ 0.5 \end{bmatrix}$$

2³ experiment on standardform

A	B	C	AB	AC	BC	ABC	levelcode	y
-	-	-	+	+	+	-	1	60
+	-	-	-	-	+	+	a	72
-	+	-	-	+	-	+	b	54
+	+	-	+	-	-	-	ab	68
-	-	+	+	-	-	+	c	52
+	-	+	-	+	-	-	ac	83
-	+	+	-	-	+	-	bc	45
+	+	+	+	+	+	+	abc	80

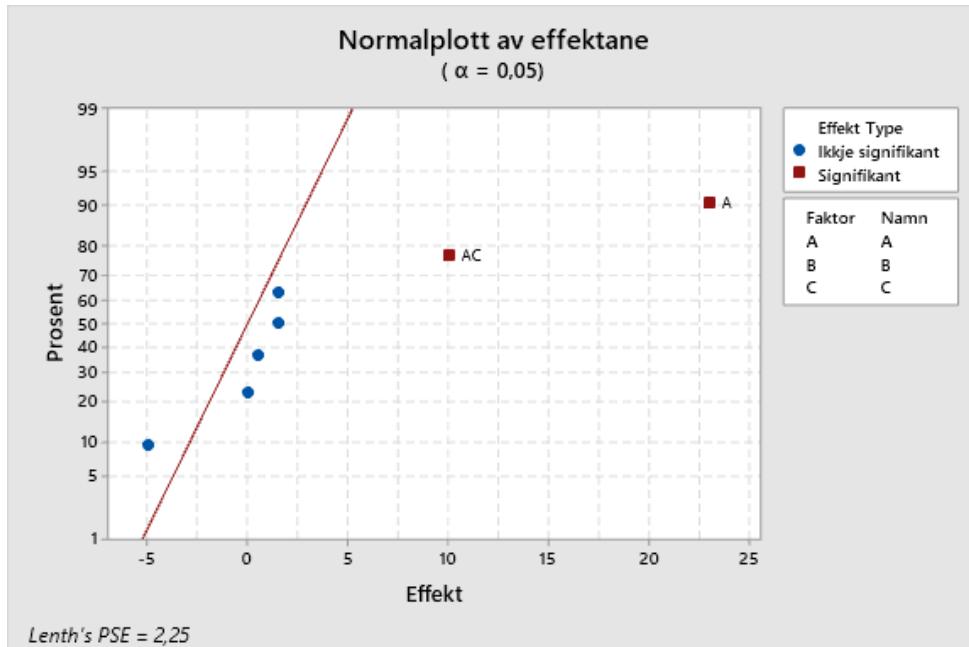
>Calulation of effects

$$\bar{y}_H - \bar{y}_L$$

Variance

$$\text{Var}(\bar{y}_H - \bar{y}_L) = \frac{4\sigma^2}{n}$$

Evaluation of significance

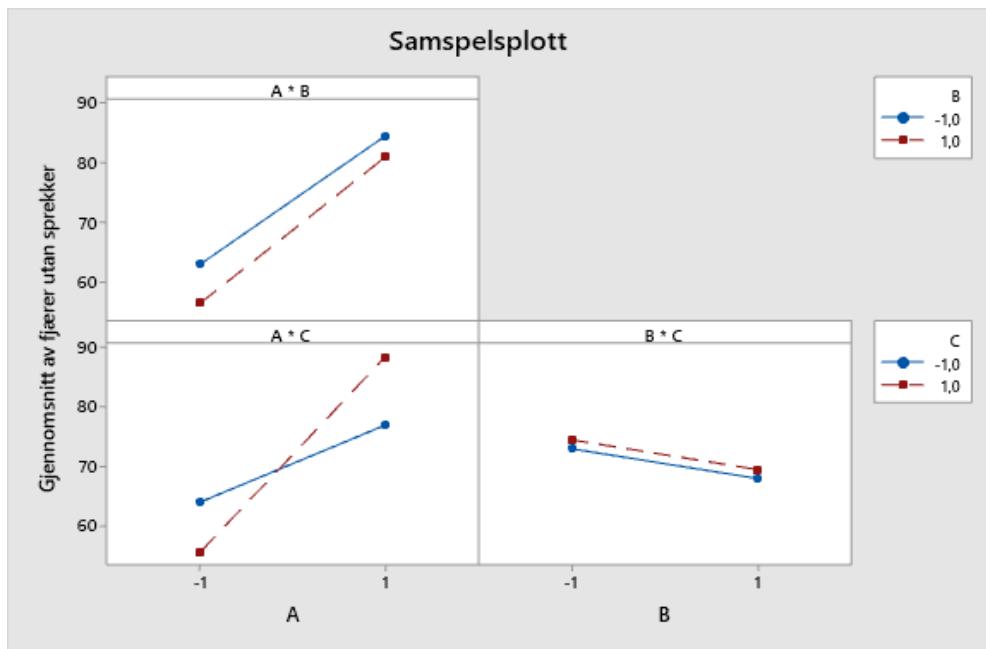
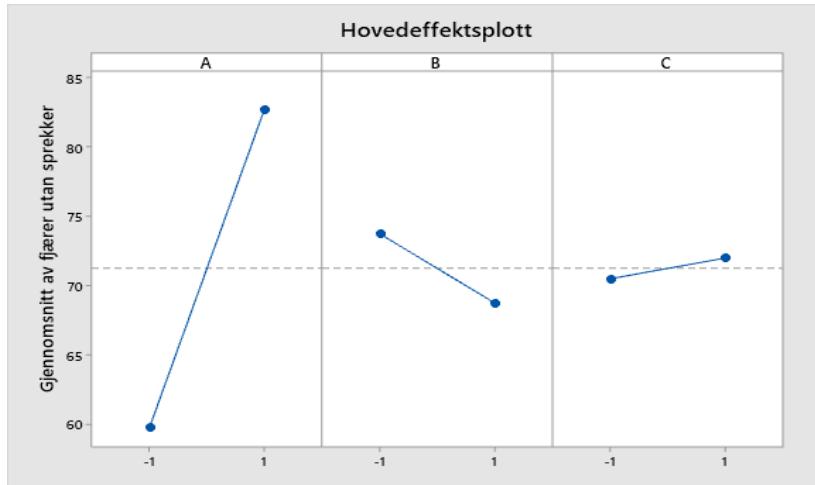


$$\sigma_{\text{effekt}}^2 \text{ is given by: } \frac{AB\hat{C}^2 + AB\hat{D}^2 + AC\hat{D}^2 + BC\hat{D}^2 + ABC\hat{D}^2}{5}$$

Replicates

Average over $\sum_{j=1}^m \frac{(Y_{ij} - \bar{Y}_i)^2}{m-1}, i = 1, 2, \dots, n$

Interpretasjon



Blocking and Fractioning

Define Block generators:

8 Blocks: $B_1 = ACE$ $B_2 = ABEF$ and $B_3 = ABCD$.

$$B_1 B_2 = ACEABEF = BCF$$

$$B_1 B_3 = ACEABCD = BDE$$

$$B_2 B_3 = ABEFABCD = CDEF$$

$$B_1 B_2 B_3 = ACEABEFABCD = ADF$$

Fractional factorials

Find defining relation from generators.

$$2^{5-2} : D=AB, E=AC$$

$$I=ABD \text{ and } I=ACE \text{ and } I^2=I=ABDACE=BCDE.$$

Hence the defining relation is $I=ABD=ACE=BCDE$.

$$2^{7-4} : D=AB, E=AC, F=BC \text{ and } G=ABC$$

$$I = ABD=ACE=BCF=ABCG$$

$$I^2=I=BCDE=ACDF=CDG=ABEF=BEG=AFG$$

$$I^3=I=DEF=ADEG=BDFG=CEFG$$

$$I^4=I=ABCDEFG$$

Analysis of variance

One-way

$$Y_{ij} = \mu + \alpha_i + \varepsilon_{ij} \begin{cases} N(0, \sigma^2), \text{ and independent} \\ i = 1, 2, \dots, k, j = 1, 2, \dots, n_i \end{cases}$$

Hypothesis

$$H_0: \alpha_1 = \alpha_2 = \dots = \alpha_k = 0 \quad H_1: \text{at least one } \alpha_i \neq 0.$$

Source	SS	DF	MS	F
Treatment	$SS_A = \sum_{i=1}^k n_i (\bar{y}_{i\cdot} - \bar{y}_{..})^2$	$k-1$	$SS_A / k-1$	$\frac{SS_A}{k-1} / \frac{SS_E}{n-k}$
Error	$SS_E = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\cdot})^2$	$n-k$	$SS_E / n-k$	
Total	$SS_T = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{..})^2$	$n-1$		

Randomized complete block design

$$Y_{ij} = \mu + \alpha_i + \beta_j + \varepsilon_{ij} \begin{cases} N(0, \sigma^2) \text{ and independent} \\ j = 1, 2, \dots, b, \quad i = 1, 2, \dots, k \end{cases}$$

Hypothesis

$$H_0: \alpha_1 = \alpha_2 = \dots = \alpha_k = 0 \quad H_1: \text{at least two are different from zero}$$

Sources	SS	DF	MS	F
Treatment	$SS_A = b \sum_{i=1}^k (\bar{y}_{i\cdot} - \bar{y}_{..})^2$	$k-1$	$SS_A / k-1$	$\frac{SS_A}{k-1} / \frac{SS_E}{b-1} \frac{k-1}{k-1}$
Block	$SS_B = k \sum_{j=1}^b (\bar{y}_{.j} - \bar{y}_{..})^2$	$b-1$	$SS_B / b-1$	
Error	$SS_E = \sum_{i=1}^k \sum_{j=1}^b (y_{ij} - \bar{y}_{i\cdot} - \bar{y}_{.j} + \bar{y}_{..})^2$	$(b-1)(k-1)$	$SS_E / b-1$	$k-1$
Total	$SS_T = \sum_{i=1}^k \sum_{j=1}^b (y_{ij} - \bar{y}_{..})^2$	$bk-1$		

Two-way analysis of variance

Model

$$Y_{ijk} = \mu + \alpha_i + \beta_j + \alpha\beta_{ij} + \varepsilon_{ijk} \quad \begin{cases} \text{N } 0, \sigma^2 \text{ and independent} \\ i=1, 2, \dots, a, \ j=1, 2, \dots, b, \ k=1, 2, \dots, n \end{cases}$$

1. $H_0: \alpha_1 = \alpha_2 = \dots = \alpha_a = 0 \quad H_1: \text{at least one is different from 0.}$
2. $H_0: \beta_1 = \beta_2 = \dots = \beta_b = 0 \quad H_1: \text{at least one is different from 0.}$
3. $H_0: \alpha\beta_{11} = \alpha\beta_{12} = \dots = \alpha\beta_{ab} = 0 \quad H_1: \text{at least one is different from 0.}$

Sources	SS	DF	MS	F
A	$SS_A = nb \sum_{i=1}^a y_{i..} - y_{...}^2$	$a-1$	$\frac{SS_A}{a-1}$	$F = \frac{MS_A}{MS_E}$
B	$SS_B = na \sum_{j=1}^b y_{.j.} - \bar{y}_{...}^2$	$b-1$	$\frac{SS_B}{b-1}$	$F = \frac{MS_B}{MS_E}$
Interaction AB	$SS_{AB} = n \sum_{i=1}^a \sum_{j=1}^b y_{ij.} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y}_{...}^2$	$a-1 \ b-1$	$\frac{SS_{AB}}{a-1 \ b-1}$	$F = \frac{MS_{AB}}{MS_E}$
Error	$SS_E = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^n y_{ijk} - \bar{y}_{ij.}^2$	$ab \ n-1$	$\frac{SS_E}{ab \ n-1}$	
Total	$SS_T = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^n y_{ijk} - \bar{y}_{...}^2$	$abn-1$		

Control Charts

Let T be a test statistics based on a sample X_1, X_2, \dots, X_n .

Control Chart: $LCL = \mu_T - k\sigma_T$, $CL = \mu_T$, $UCL = \mu_T + k\sigma_T$

Reduce the probability of false alarm: increase k.

Increase the probability of detecting out of control: increase n.

\bar{X} -R Chart

k samples: Test statistics: $\bar{X}_1, \dots, \bar{X}_k, R_1, \dots, R_k$.

$$\bar{\bar{X}} = \frac{\sum_{i=1}^k \bar{X}_i}{k}, \quad \bar{R} = \frac{\sum_{i=1}^k R_i}{k}.$$

$$\bar{X}\text{-bar: } LCL = \bar{\bar{X}} + \frac{3\bar{R}}{d_2\sqrt{n}}, \quad CL = \bar{\bar{X}}, \quad UCL = \bar{\bar{X}} + \frac{3\bar{R}}{d_2\sqrt{n}}.$$

$$R \text{ Chart: } LCL = \bar{R} - \frac{3\bar{R}d_3}{d_2}, \quad CL = \bar{R}, \quad UCL = \bar{R} + \frac{3\bar{R}d_3}{d_2}$$

p-charts

$$LCL = \bar{p} - 3\sqrt{\frac{\bar{p}(1-\bar{p})}{n}}, \quad CL = \bar{p}, \quad UCL = \bar{p} + 3\sqrt{\frac{\bar{p}(1-\bar{p})}{n}}$$

C-Charts

$$LCL = \hat{\lambda} - 3\sqrt{\hat{\lambda}}, \quad CL = \hat{\lambda}, \quad UCL = \hat{\lambda} + 3\sqrt{\hat{\lambda}}$$

χ^2 - tests

GOF

$$\chi^2 = \sum_{i=1}^k \frac{(O_i - e_i)^2}{e_i} \approx \chi^2(k-1)$$

$$e_i = np_i$$

$r \times c$ Contingency Tables

Test for independence : $H_0 : p_{ij} = p_i \cdot p_j \quad \forall i, j$

Test for equal proportions in populations:

$$H_0 : p_{ij} = p_i \quad \forall j \text{ and for all } i.$$

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - e_{ij})^2}{e_{ij}} \approx \chi^2((r-1)(c-1))$$

$$e_{ij} = \frac{n_i \cdot n_j}{n}$$

Nonparametric tests

Sign test (continuous distribution)

X = Number of + (or -) for $x_i - \tilde{\mu}_0$, $i=1,2,\dots,n$

Test statistics: $X \sim B(n, 0.5)$

Wilcoxon 1 sample test (continuous and symmetric distribution)

Rank $|x_i - \tilde{\mu}_0|$, $i=1,2,\dots,n$ in ascending order with ranks given by $1,2,\dots,n$.

Test statistics: $W_+ = \sum_{i:x_i - \tilde{\mu}_0 > 0} R_i$ and $W_- = \sum_{i:x_i - \tilde{\mu}_0 < 0} R_i$

Wilcoxon 2 sample test (continuous and identical distributions except for location)

Rank all observations in ascending order, in total $n_1 + n_2$

Test statistics: $U_1 = \sum_{x_i \in \text{sample 1}} R_i - \frac{n_1(n_1 + 1)}{2}$ and

$U_2 = \sum_{x_i \in \text{sample 2}} R_i - \frac{n_2(n_2 + 1)}{2}$

Kruskal Wallis test for number of samples, $k > 2$ (continuous and identical distributions except for location)

Rank all observations in ascending order, in total $n_1 + n_2 + \dots + n_k = n$

Compute average rank, R_i , for each sample.

$$\text{Test statistics: } H = \frac{12}{n(n+1)} \sum_{i=1}^k \frac{R_i^2}{n_i} - 3(n+1) \approx \chi^2(k-1)$$